

FIGURE 1.1: A “particle” constrained to move in one dimension under the influence of a specified force.

Here i is the square root of -1 , and \hbar is Planck’s constant—or rather, his *original* constant (h) divided by 2π :

$$\hbar = \frac{h}{2\pi} = 1.054572 \times 10^{-34} \text{ J s.} \quad [1.2]$$

The Schrödinger equation plays a role logically analogous to Newton’s second law: Given suitable initial conditions (typically, $\Psi(x, 0)$), the Schrödinger equation determines $\Psi(x, t)$ for all future time, just as, in classical mechanics, Newton’s law determines $x(t)$ for all future time.²

1.2 THE STATISTICAL INTERPRETATION

But what exactly *is* this “wave function,” and what does it do for you once you’ve *got* it? After all, a particle, by its nature, is localized at a point, whereas the wave function (as its name suggests) is spread out in space (it’s a function of x , for any given time t). How can such an object represent the state of a *particle*? The answer is provided by Born’s **statistical interpretation** of the wave function, which says that $|\Psi(x, t)|^2$ gives the *probability* of finding the particle at point x , at time t —or, more precisely,³

$$\int_a^b |\Psi(x, t)|^2 dx = \left\{ \begin{array}{l} \text{probability of finding the particle} \\ \text{between } a \text{ and } b, \text{ at time } t. \end{array} \right\} \quad [1.3]$$

Probability is the *area* under the graph of $|\Psi|^2$. For the wave function in Figure 1.2, you would be quite likely to find the particle in the vicinity of point A , where $|\Psi|^2$ is large, and relatively *unlikely* to find it near point B .

²For a delightful first-hand account of the origins of the Schrödinger equation see the article by Felix Bloch in *Physics Today*, December 1976.

³The wave function itself is complex, but $|\Psi|^2 = \Psi^* \Psi$ (where Ψ^* is the complex conjugate of Ψ) is real and nonnegative—as a probability, of course, *must* be.

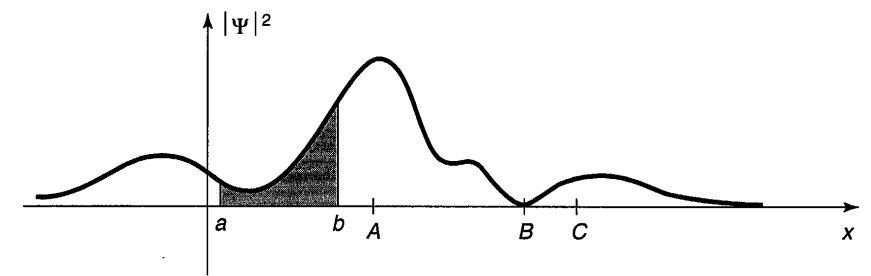


FIGURE 1.2: A typical wave function. The shaded area represents the probability of finding the particle between a and b . The particle would be relatively likely to be found near A , and unlikely to be found near B .

The statistical interpretation introduces a kind of **indeterminacy** into quantum mechanics, for even if you know everything the theory has to tell you about the particle (to wit: its wave function), still you cannot predict with certainty the outcome of a simple experiment to measure its position—all quantum mechanics has to offer is *statistical* information about the *possible* results. This indeterminacy has been profoundly disturbing to physicists and philosophers alike, and it is natural to wonder whether it is a fact of nature, or a defect in the theory.

Suppose I *do* measure the position of the particle, and I find it to be at point C .⁴ *Question:* Where was the particle just *before* I made the measurement? There are three plausible answers to this question, and they serve to characterize the main schools of thought regarding quantum indeterminacy:

1. The realist position: *The particle was at C.* This certainly seems like a sensible response, and it is the one Einstein advocated. Note, however, that if this is true then quantum mechanics is an *incomplete* theory, since the particle *really* was at C , and yet quantum mechanics was unable to tell us so. To the realist, indeterminacy is not a fact of nature, but a reflection of our ignorance. As d’Espagnat put it, “the position of the particle was never indeterminate, but was merely unknown to the experimenter.”⁵ Evidently Ψ is not the whole story—some additional information (known as a **hidden variable**) is needed to provide a complete description of the particle.

2. The orthodox position: *The particle wasn’t really anywhere.* It was the act of measurement that forced the particle to “take a stand” (though how and why it decided on the point C we dare not ask). Jordan said it most starkly: “Observations not only *disturb* what is to be measured, they *produce* it . . . We *compel* (the

⁴Of course, no measuring instrument is perfectly precise; what I *mean* is that the particle was found *in the vicinity* of C , to within the tolerance of the equipment.

⁵Bernard d’Espagnat, “The Quantum Theory and Reality” (*Scientific American*, November 1979, p. 165).

particle) to assume a definite position.”⁶ This view (the so-called **Copenhagen interpretation**), is associated with Bohr and his followers. Among physicists it has always been the most widely accepted position. Note, however, that if it is correct there is something very peculiar about the act of measurement—something that over half a century of debate has done precious little to illuminate.

3. The **agnostic** position: *Refuse to answer*. This is not quite as silly as it sounds—after all, what sense can there be in making assertions about the status of a particle *before* a measurement, when the only way of knowing whether you were right is precisely to conduct a measurement, in which case what you get is no longer “before the measurement?” It is metaphysics (in the pejorative sense of the word) to worry about something that cannot, by its nature, be tested. Pauli said: “One should no more rack one’s brain about the problem of whether something one cannot know anything about exists all the same, than about the ancient question of how many angels are able to sit on the point of a needle.”⁷ For decades this was the “fall-back” position of most physicists: They’d try to sell you the orthodox answer, but if you were persistent they’d retreat to the agnostic response, and terminate the conversation.

Until fairly recently, all three positions (realist, orthodox, and agnostic) had their partisans. But in 1964 John Bell astonished the physics community by showing that it makes an *observable* difference whether the particle had a precise (though unknown) position prior to the measurement, or not. Bell’s discovery effectively eliminated agnosticism as a viable option, and made it an *experimental* question whether 1 or 2 is the correct choice. I’ll return to this story at the end of the book, when you will be in a better position to appreciate Bell’s argument; for now, suffice it to say that the experiments have decisively confirmed the orthodox interpretation.⁸ A particle simply does not *have* a precise position prior to measurement, any more than the ripples on a pond do; it is the measurement process that insists on one particular number, and thereby in a sense *creates* the specific result, limited only by the statistical weighting imposed by the wave function.

What if I made a *second* measurement, *immediately* after the first? Would I get C again, or does the act of measurement cough up some completely new number each time? On this question everyone is in agreement: A repeated measurement (on the same particle) must return the same value. Indeed, it would be tough to prove that the particle was really found at C in the first instance, if this could not be confirmed by immediate repetition of the measurement. How does the orthodox

⁶Quoted in a lovely article by N. David Mermin, “Is the moon there when nobody looks?” (Physics Today, April 1985, p. 38).

⁷Quoted by Mermin (footnote 6), p. 40.

⁸This statement is a little too strong: There remain a few theoretical and experimental loopholes, some of which I shall discuss in the Afterword. There exist viable nonlocal hidden variable theories (notably David Bohm’s), and other formulations (such as the **many worlds** interpretation) that do not fit cleanly into any of my three categories. But I think it is wise, at least from a pedagogical point of view, to adopt a clear and coherent platform at this stage, and worry about the alternatives later.

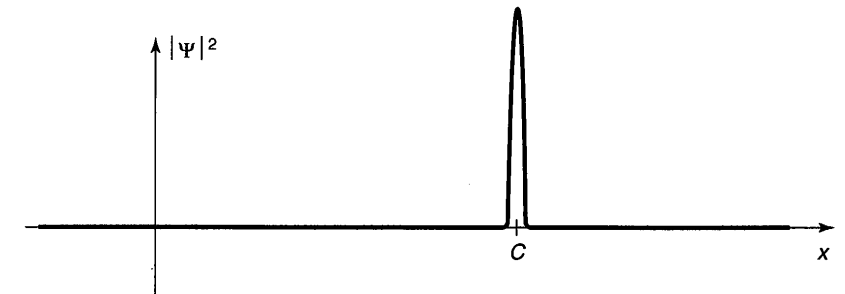


FIGURE 1.3: Collapse of the wave function: graph of $|\Psi|^2$ immediately after a measurement has found the particle at point C .

interpretation account for the fact that the second measurement is bound to yield the value C ? Evidently the first measurement radically alters the wave function, so that it is now sharply peaked about C (Figure 1.3). We say that the wave function **collapses**, upon measurement, to a spike at the point C (it soon spreads out again, in accordance with the Schrödinger equation, so the second measurement must be made quickly). There are, then, two entirely distinct kinds of physical processes: “ordinary” ones, in which the wave function evolves in a leisurely fashion under the Schrödinger equation, and “measurements,” in which Ψ suddenly and discontinuously collapses.⁹

1.3 PROBABILITY

1.3.1 Discrete Variables

Because of the statistical interpretation, probability plays a central role in quantum mechanics, so I digress now for a brief discussion of probability theory. It is mainly a question of introducing some notation and terminology, and I shall do it in the context of a simple example.

Imagine a room containing fourteen people, whose ages are as follows:

- one person aged 14,
- one person aged 15,
- three people aged 16,

⁹The role of measurement in quantum mechanics is so critical and so bizarre that you may well be wondering what precisely *constitutes* a measurement. Does it have to do with the interaction between a microscopic (quantum) system and a macroscopic (classical) measuring apparatus (as Bohr insisted), or is it characterized by the leaving of a permanent “record” (as Heisenberg claimed), or does it involve the intervention of a conscious “observer” (as Wigner proposed)? I’ll return to this thorny issue in the Afterword; for the moment let’s take the naive view: A measurement is the kind of thing that a scientist does in the laboratory, with rulers, stopwatches, Geiger counters, and so on.

AFTERWORD

Now that you have (I hope) a sound understanding of what quantum mechanics *says*, I would like to return to the question of what it *means*—continuing the story begun in Section 1.2. The source of the problem is the indeterminacy associated with the statistical interpretation of the wave function. For Ψ (or, more generally, the *quantum state*—it could be a spinor, for example) does not uniquely determine the outcome of a measurement; all it provides is the statistical distribution of possible results. This raises a profound question: Did the physical system “actually have” the attribute in question *prior* to the measurement (the so-called **realist** viewpoint), or did the act of measurement itself “create” the property, limited only by the statistical constraint imposed by the wave function (the **orthodox** position)—or can we duck the question entirely, on the grounds that it is “metaphysical” (the **agnostic** response)?

According to the realist, quantum mechanics is an *incomplete* theory, for even if you know *everything quantum mechanics has to tell you* about the system (to wit: its wave function), still you cannot determine all of its features. Evidently there is some *other* information, external to quantum mechanics, which (together with Ψ) is required for a complete description of physical reality.

The orthodox position raises even more disturbing problems, for if the act of measurement forces the system to “take a stand,” helping to *create* an attribute that was not there previously,¹ then there is something very peculiar about the

¹This may be *strange*, but it is not *mystical*, as some popularizers would like to suggest. The so-called **wave-particle duality**, which Niels Bohr elevated to the status of a cosmic principle (**complementarity**), makes electrons sound like unpredictable adolescents, who sometimes behave like

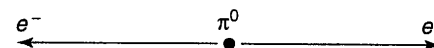


FIGURE 12.1: Bohm’s version of the EPR experiment: A π^0 at rest decays into an electron-positron pair.

measurement process. Moreover, in order to account for the fact that an immediately repeated measurement yields the same result, we are forced to assume that the act of measurement **collapses** the wave function, in a manner that is difficult, at best, to reconcile with the normal evolution prescribed by the Schrödinger equation.

In light of this, it is no wonder that generations of physicists retreated to the agnostic position, and advised their students not to waste time worrying about the conceptual foundations of the theory.

12.1 THE EPR PARADOX

In 1935, Einstein, Podolsky, and Rosen² published the famous **EPR paradox**, which was designed to prove (on purely theoretical grounds) that the realist position is the only sustainable one. I’ll describe a simplified version of the EPR paradox, introduced by David Bohm. Consider the decay of the neutral pi meson into an electron and a positron:

$$\pi^0 \rightarrow e^- + e^+.$$

Assuming the pion was at rest, the electron and positron fly off in opposite directions (Figure 12.1). Now, the pion has spin zero, so conservation of angular momentum requires that the electron and positron are in the singlet configuration:

$$\frac{1}{\sqrt{2}}(\uparrow\downarrow + \downarrow\uparrow). \quad [12.1]$$

If the electron is found to have spin up, the positron must have spin down, and vice versa. Quantum mechanics can’t tell you *which* combination you’ll get, in any particular pion decay, but it does say that the measurements will be *correlated*, and you’ll get each combination half the time (on average). Now suppose

adults, and sometimes, for no particular reason, like children. I prefer to avoid such language. When I say that a particle does not have a particular attribute before its measurement, I have in mind, for example, an electron in the spin state $\chi = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$; a measurement of the x -component of its angular momentum could return the value $\hbar/2$, or (with equal probability) the value $-\hbar/2$, but until the measurement is made it simply *does not have* a well-defined value of S_x .

²A. Einstein, B. Podolsky, and N. Rosen, *Phys. Rev.* **47**, 777 (1935).

we let the electron and positron fly way off—10 meters, in a practical experiment, or, in principle, 10 light years—and then you measure the spin of the electron. Say you get spin up. Immediately you know that someone 20 meters (or 20 light years) away will get spin down, if that person examines the positron.

To the realist, there's nothing surprising in this—the electron *really had* spin up (and the positron spin down) from the moment they were created ... it's just that quantum mechanics didn't know about it. But the “orthodox” view holds that neither particle had either spin up *or* spin down until the act of measurement intervened: Your measurement of the electron collapsed the wave function, and instantaneously “produced” the spin of the positron 20 meters (or 20 light years) away. Einstein, Podolsky, and Rosen considered such “spooky action-at-a-distance” (Einstein's words) preposterous. They concluded that the orthodox position is untenable; the electron and positron must have had well-defined spins all along, whether quantum mechanics can calculate them or not.

The fundamental assumption on which the EPR argument rests is that no influence can propagate faster than the speed of light. We call this the principle of **locality**. You might be tempted to propose that the collapse of the wave function is *not* instantaneous, but “travels” at some finite velocity. However, this would lead to violations of angular momentum conservation, for if we measured the spin of the positron before the news of the collapse had reached it, there would be a fifty-fifty probability of finding *both* particles with spin up. Whatever one might think of such a theory in the abstract, the experiments are unequivocal: No such violation occurs—the (anti-)correlation of the spins is perfect. Evidently the collapse of the wave function—whatever its ontological status—is instantaneous.

Problem 12.1 Entangled states. The singlet spin configuration (Equation 12.1) is the classic example of an *entangled state*—a two-particle state that cannot be expressed as the product of two one-particle states, and for which, therefore, one cannot really speak of “the state” of either particle separately. You might wonder whether this is somehow an artifact of bad notation—maybe some linear combination of the one-particle states would disentangle the system. Prove the following theorem:

Consider a two-level system, $|\phi_a\rangle$ and $|\phi_b\rangle$, with $\langle\phi_i|\phi_j\rangle = \delta_{ij}$. (For example, $|\phi_a\rangle$ might represent spin up and $|\phi_b\rangle$ spin down.) The two-particle state

$$\alpha|\phi_a(1)\rangle|\phi_b(2)\rangle + \beta|\phi_b(1)\rangle|\phi_a(2)\rangle$$

(with $\alpha \neq 0$ and $\beta \neq 0$) *cannot* be expressed as a product

$$|\psi_r(1)\rangle|\psi_s(2)\rangle,$$

for *any* one-particle states $|\psi_r\rangle$ and $|\psi_s\rangle$.

Hint: Write $|\psi_r\rangle$ and $|\psi_s\rangle$ as linear combinations of $|\phi_a\rangle$ and $|\phi_b\rangle$.

12.2 BELL'S THEOREM

Einstein, Podolsky, and Rosen did not doubt that quantum mechanics is *correct*, as far as it goes; they only claimed that it is an *incomplete* description of physical reality: The wave function is not the whole story—some *other* quantity, λ , is needed, in addition to Ψ , to characterize the state of a system fully. We call λ the “hidden variable” because, at this stage, we have no idea how to calculate or measure it.³ Over the years, a number of hidden variable theories have been proposed, to supplement quantum mechanics;⁴ they tend to be cumbersome and implausible, but never mind—until 1964 the program seemed eminently worth pursuing. But in that year J. S. Bell proved that *any* local hidden variable theory is *incompatible* with quantum mechanics.⁵

Bell suggested a generalization of the EPR/Bohm experiment: Instead of orienting the electron and positron detectors along the *same* direction, he allowed them to be rotated independently. The first measures the component of the electron spin in the direction of a unit vector \mathbf{a} , and the second measures the spin of the positron along the direction \mathbf{b} (Figure 12.2). For simplicity, let's record the spins in units of $\hbar/2$; then each detector registers the value $+1$ (for spin up) or -1 (spin down), along the direction in question. A table of results, for many π^0 decays, might look like this:

electron	positron	product
+1	-1	-1
+1	+1	+1
-1	+1	-1
+1	-1	-1
-1	-1	+1
⋮	⋮	⋮

³The hidden variable could be a single number, or it could be a whole *collection* of numbers; perhaps λ is to be calculated in some future theory, or maybe it is for some reason of principle incalculable. It hardly matters. All I am asserting is that there must be *something*—if only a *list* of the outcomes of every possible experiment—associated with the system prior to a measurement.

⁴D. Bohm, *Phys. Rev.* **85**, 166, 180 (1952).

⁵Bell's original paper (*Physics* **1**, 195 (1964)) is a gem: brief, accessible, and beautifully written.

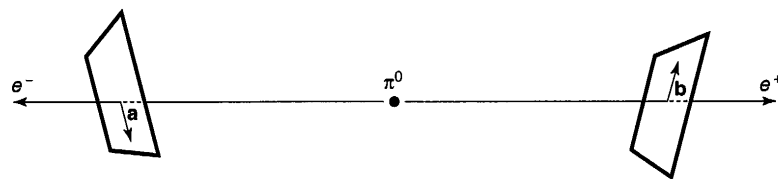


FIGURE 12.2: Bell's version of the EPR-Bohm experiment: detectors independently oriented in directions \mathbf{a} and \mathbf{b} .

Bell proposed to calculate the *average* value of the *product* of the spins, for a given set of detector orientations. Call this average $P(\mathbf{a}, \mathbf{b})$. If the detectors are parallel ($\mathbf{b} = \mathbf{a}$), we recover the original EPRB configuration; in this case one is spin up and the other spin down, so the product is always -1 , and hence so too is the average:

$$P(\mathbf{a}, \mathbf{a}) = -1. \quad [12.2]$$

By the same token, if they are *anti*-parallel ($\mathbf{b} = -\mathbf{a}$), then every product is $+1$, so

$$P(\mathbf{a}, -\mathbf{a}) = +1. \quad [12.3]$$

For arbitrary orientations, quantum mechanics predicts

$$P(\mathbf{a}, \mathbf{b}) = -\mathbf{a} \cdot \mathbf{b} \quad [12.4]$$

(see Problem 4.50). What Bell discovered is that *this result is incompatible with any local hidden variable theory*.

The argument is stunningly simple. Suppose that the “complete” state of the electron/positron system is characterized by the hidden variable(s) λ (λ varies, in some way that we neither understand nor control, from one pion decay to the next). Suppose further that the outcome of the *electron* measurement is independent of the orientation (\mathbf{b}) of the *positron* detector—which may, after all, be chosen by the experimenter at the positron end just before the electron measurement is made, and hence far too late for any subluminal message to get back to the electron detector. (This is the locality assumption.) Then there exists some function $A(\mathbf{a}, \lambda)$ which gives the result of an electron measurement, and some other function $B(\mathbf{b}, \lambda)$ for the positron measurement. These functions can only take on the values ± 1 :⁶

$$A(\mathbf{a}, \lambda) = \pm 1; \quad B(\mathbf{b}, \lambda) = \pm 1. \quad [12.5]$$

⁶This already concedes far more than a *classical* determinist would be prepared to allow, for it abandons any notion that the particles could have well-defined angular momentum vectors with simultaneously determinate components. But never mind—the point of Bell's argument is to demonstrate that quantum mechanics is incompatible with *any* local deterministic theory—even one that bends over backwards to be accommodating.

When the detectors are aligned, the results are perfectly (anti)-correlated:

$$A(\mathbf{a}, \lambda) = -B(\mathbf{a}, \lambda), \quad [12.6]$$

for all λ .

Now, the average of the product of the measurements is

$$P(\mathbf{a}, \mathbf{b}) = \int \rho(\lambda) A(\mathbf{a}, \lambda) B(\mathbf{b}, \lambda) d\lambda, \quad [12.7]$$

where $\rho(\lambda)$ is the probability density for the hidden variable. (Like any probability density, it is nonnegative, and satisfies the normalization condition $\int \rho(\lambda) d\lambda = 1$, but beyond this we make no assumptions about $\rho(\lambda)$; different hidden variable theories would presumably deliver quite different expressions for ρ .) In view of Equation 12.6, we can eliminate B :

$$P(\mathbf{a}, \mathbf{b}) = - \int \rho(\lambda) A(\mathbf{a}, \lambda) A(\mathbf{b}, \lambda) d\lambda, \quad [12.8]$$

If \mathbf{c} is any *other* unit vector,

$$P(\mathbf{a}, \mathbf{b}) - P(\mathbf{a}, \mathbf{c}) = - \int \rho(\lambda) [A(\mathbf{a}, \lambda) A(\mathbf{b}, \lambda) - A(\mathbf{a}, \lambda) A(\mathbf{c}, \lambda)] d\lambda. \quad [12.9]$$

Or, since $[A(\mathbf{b}, \lambda)]^2 = 1$:

$$P(\mathbf{a}, \mathbf{b}) - P(\mathbf{a}, \mathbf{c}) = - \int \rho(\lambda) [1 - A(\mathbf{b}, \lambda) A(\mathbf{c}, \lambda)] A(\mathbf{a}, \lambda) A(\mathbf{b}, \lambda) d\lambda. \quad [12.10]$$

But it follows from Equation 12.5 that $-1 \leq [A(\mathbf{a}, \lambda) A(\mathbf{b}, \lambda)] \leq +1$; moreover $\rho(\lambda) [1 - A(\mathbf{b}, \lambda) A(\mathbf{c}, \lambda)] \geq 0$, so

$$|P(\mathbf{a}, \mathbf{b}) - P(\mathbf{a}, \mathbf{c})| \leq \int \rho(\lambda) [1 - A(\mathbf{b}, \lambda) A(\mathbf{c}, \lambda)] d\lambda, \quad [12.11]$$

or, more simply:

$$|P(\mathbf{a}, \mathbf{b}) - P(\mathbf{a}, \mathbf{c})| \leq 1 + P(\mathbf{b}, \mathbf{c}). \quad [12.12]$$

This is the famous **Bell inequality**. It holds for *any* local hidden variable theory (subject only to the minimal requirements of Equations 12.5 and 12.6), for we have made no assumptions whatever as to the nature or number of the hidden variable(s), or their distribution (ρ).

But it is easy to show that the quantum mechanical prediction (Equation 12.4) is incompatible with Bell's inequality. For example, suppose all three vectors lie in

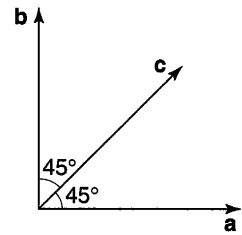


FIGURE 12.3: An orientation of the detectors that demonstrates quantum violations of Bell's inequality.

a plane, and c makes a 45° angle with a and b (Figure 12.3); in this case quantum mechanics says

$$P(\mathbf{a}, \mathbf{b}) = 0, \quad P(\mathbf{a}, \mathbf{c}) = P(\mathbf{b}, \mathbf{c}) = -0.707,$$

which is patently inconsistent with Bell's inequality:

$$0.707 \not\leq 1 - 0.707 = 0.293.$$

With Bell's modification, then, the EPR paradox proves something far more radical than its authors imagined: If they are right, then not only is quantum mechanics *incomplete*, it is downright *wrong*. On the other hand, if quantum mechanics is right, then *no* hidden variable theory is going to rescue us from the nonlocality Einstein considered so preposterous. Moreover, we are provided with a very simple experiment to settle the issue once and for all.

Many experiments to test Bell's inequality were performed in the '60's and '70's, culminating in the work of Aspect, Grangier, and Roger.⁷ The details do not concern us here (they actually used two-photon atomic transitions, not pion decays). To exclude the remote possibility that the positron detector might somehow "sense" the orientation of the electron detector, both orientations were set quasi-randomly *after* the photons were already in flight. The results were in excellent agreement with the predictions of quantum mechanics, and clearly incompatible with Bell's inequality.⁸

Ironically, the experimental confirmation of quantum mechanics came as something of a shock to the scientific community. But not because it spelled the demise of "realism"—most physicists had long since adjusted to this (and for those who could not, there remained the possibility of *nonlocal* hidden variable theories,

⁷A. Aspect, P. Grangier, and G. Roger, *Phys. Rev. Lett.* **49**, 91 (1982). For more recent experiments see G. Weihs *et al.*, *Phys. Rev. Lett.* **81**, 5039 (1998).

⁸Bell's theorem involves *averages* and it is conceivable that an apparatus such as Aspect's contains some secret bias which selects out a nonrepresentative sample, thus distorting the average. In 1989, an improved version of Bell's theorem was proposed, in which a *single measurement* suffices to distinguish between the quantum prediction and that of any local hidden variable theory. See D. Greenberger, M. Horne, A. Shimony, and A. Zeilinger, *Am. J. Phys.* **58**, 1131 (1990) and N. David Mermin, *Am. J. Phys.* **58**, 731 (1990).

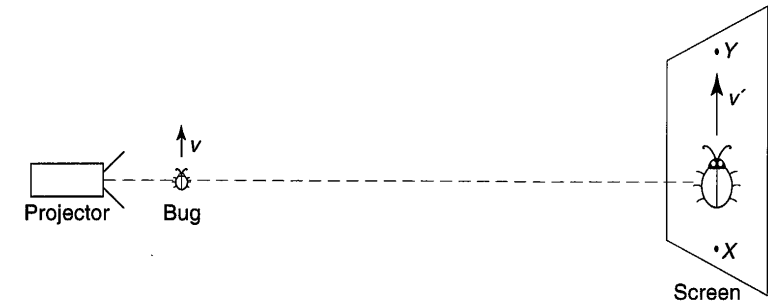


FIGURE 12.4: The shadow of the bug moves across the screen at a velocity v' greater than c , provided the screen is far enough away.

to which Bell's theorem does not apply⁹). The real shock was the demonstration that *nature itself is fundamentally nonlocal*. Nonlocality, in the form of the instantaneous collapse of the wave function (and for that matter also in the symmetrization requirement for identical particles) had always been a feature of the orthodox interpretation, but before Aspect's experiment it was possible to hope that quantum nonlocality was somehow a nonphysical artifact of the formalism, with no detectable consequences. That hope can no longer be sustained, and we are obliged to reexamine our objection to instantaneous action-at-a-distance.

Why *are* physicists so squeamish about superluminal influences? After all, there are many things that travel faster than light. If a bug flies across the beam of a movie projector, the speed of its shadow is proportional to the distance to the screen; in principle, that distance can be as large as you like, and hence the *shadow* can move at arbitrarily high velocity (Figure 12.4). However, the shadow does not carry any *energy*; nor can it transmit a *message* from one point to another on the screen. A person at point X cannot *cause anything to happen* at point Y by manipulating the passing shadow.

On the other hand, a *causal* influence that propagated faster than light would carry unacceptable implications. For according to special relativity there exist inertial frames in which such a signal propagates *backward in time*—the effect preceding the cause—and this leads to inescapable logical anomalies. (You could, for example, arrange to kill your infant grandfather. Not a good idea!) The question is, are the superluminal influences predicted by quantum mechanics and detected by

⁹It is a curious twist of fate that the EPR paradox, which *assumed* locality in order to *prove* realism, led finally to the demise of locality and left the issue of realism undecided—the outcome (as Mermin put it) Einstein would have liked *least*. Most physicists today consider that if they can't have *local* realism, there's not much point in realism at *all*, and for this reason nonlocal hidden variable theories occupy a rather peripheral place. Still, some authors—notably Bell himself, in *Speakable and Unspeakable in Quantum Mechanics* (Cambridge University Press, 1987)—argue that such theories offer the best hope of bridging the conceptual gap between the measured system and the measuring apparatus, and for supplying an intelligible mechanism for the collapse of the wave function.

Aspect *causal*, in this sense, or are they somehow ethereal enough (like the motion of the shadow) to escape the philosophical objection?

Well, let's consider Bell's experiment. Does the measurement of the electron *influence* the outcome of the positron measurement? Assuredly it *does*—otherwise we cannot account for the correlation of the data. But does the measurement of the electron *cause* a particular outcome for the positron? Not in any ordinary sense of the word. There is no way the person manning the electron detector could use his measurement to send a signal to the person at the positron detector, since he does not control the outcome of his own measurement (he cannot *make* a given electron come out spin up, any more than the person at *X* can affect the passing shadow of the bug). It is true that he can decide *whether to make a measurement at all*, but the positron monitor, having immediate access only to data at his end of the line, cannot tell whether the electron was measured or not, for the lists of data compiled at the two ends, considered separately, are completely random. It is only when we *compare* the two lists later that we discover the remarkable correlations. In another reference frame the positron measurements occur *before* the electron measurements, and yet this leads to no logical paradox—the observed correlation is entirely symmetrical in its treatment, and it is a matter of indifference whether we say the observation of the electron influenced the measurement of the positron, or the other way around. This is a wonderfully delicate kind of influence whose only manifestation is a subtle correlation between two lists of otherwise random data.

We are led, then, to distinguish two types of influence: the “causal” variety, which produce actual changes in some physical property of the receiver, detectable by measurements on that subsystem alone, and an “ethereal” kind, which do not transmit energy or information, and for which the only evidence is a correlation in the data taken on the two separate subsystems—a correlation which by its nature cannot be detected by examining either list alone. Causal influences *cannot* propagate faster than light, but there is no compelling reason why ethereal ones should not. The influences associated with the collapse of the wave function are of the latter type, and the fact that they “travel” faster than light may be surprising, but it is not, after all, catastrophic.¹⁰

12.3 THE NO-CLONE THEOREM

Quantum measurements are typically **destructive**, in the sense that they alter the state of the system measured. This is how the uncertainty principle is enforced in the laboratory. You might wonder why we don't just make a bunch of identical copies (**clones**) of the original state, and measure *them*, leaving the system itself

¹⁰An enormous amount has been written about Bell's theorem. My favorite is an inspired essay by David Mermin in *Physics Today* (April 1985, page 38). An extensive bibliography will be found in L. E. Ballentine, *Am. J. Phys.* **55**, 785 (1987).

unscathed. It can't be done. Indeed, if you could build a cloning device (a “quantum Xerox machine”), quantum mechanics would be out the window.

For example, it would then be possible to send superluminal messages using the EPRB experiment. Say the message to be transmitted, from the operator of the positron detector to the operator of the electron detector, is either “yes” or “no.” If the message is “yes,” the sender measures S_z (of the positron). Never mind what result she gets—all that matters is that she makes the measurement, for this means that the electron is now in the definite state \uparrow or \downarrow (never mind which). The receiver immediately makes a million clones of the electron, and measures S_z on each of them. If they all yield the same answer (never mind *which* answer), we can be pretty sure that the electron *was* in fact measured, so the message is “yes.” If half of them are spin up, and half spin down, then the electron was definitely *not* measured, and the message is “no.”

But you *can't* make a quantum Xerox machine, as Wootters, Zurek, and Dieks proved in 1982.¹¹ Schematically, we want the machine to take as input a particle in state $|\psi\rangle$ (the one to be copied), plus a second particle in state $|X\rangle$ (the “blank sheet of paper”), and spit out *two* particles in the state $|\psi\rangle$ (original plus copy):

$$|\psi\rangle|X\rangle \rightarrow |\psi\rangle|\psi\rangle. \quad [12.13]$$

Suppose we have made a device that successfully clones the state $|\psi_1\rangle$:

$$|\psi_1\rangle|X\rangle \rightarrow |\psi_1\rangle|\psi_1\rangle, \quad [12.14]$$

and also works for state $|\psi_2\rangle$:

$$|\psi_2\rangle|X\rangle \rightarrow |\psi_2\rangle|\psi_2\rangle \quad [12.15]$$

($|\psi_1\rangle$ and $|\psi_2\rangle$ might be spin up and spin down, for example, if the particle is an electron). So far, so good. But what happens when we feed in a linear combination $|\psi\rangle = \alpha|\psi_1\rangle + \beta|\psi_2\rangle$? Evidently we get¹²

$$|\psi\rangle|X\rangle \rightarrow \alpha|\psi_1\rangle|\psi_1\rangle + \beta|\psi_2\rangle|\psi_2\rangle, \quad [12.16]$$

which is not at all what we wanted—what we *wanted* was

$$\begin{aligned} |\psi\rangle|X\rangle \rightarrow |\psi\rangle|\psi\rangle &= [\alpha|\psi_1\rangle + \beta|\psi_2\rangle][\alpha|\psi_1\rangle + \beta|\psi_2\rangle] \\ &= \alpha^2|\psi_1\rangle|\psi_1\rangle + \beta^2|\psi_2\rangle|\psi_2\rangle + \alpha\beta[|\psi_1\rangle|\psi_2\rangle + |\psi_2\rangle|\psi_1\rangle]. \end{aligned} \quad [12.17]$$

¹¹W. K. Wootters and W. H. Zurek, *Nature* **299**, 802 (1982); D. Dieks, *Phys. Lett. A* **92**, 271 (1982).

¹²This assumes that the device acts *linearly* on the state $|\psi\rangle$, as it must, since the time-dependent Schrödinger equation (which presumably governs the process) is linear.

You can make a machine to clone spin-up electrons and spin-down electrons, but it's going to fail for any nontrivial linear combinations. It's as though you bought a Xerox machine that copies vertical lines perfectly, and also horizontal lines, but completely distorts diagonals.

12.4 SCHRÖDINGER'S CAT

The measurement process plays a mischievous role in quantum mechanics: It is here that indeterminacy, nonlocality, the collapse of the wave function, and all the attendant conceptual difficulties arise. Absent measurement, the wave function evolves in a leisurely and deterministic way, according to the Schrödinger equation, and quantum mechanics looks like a rather ordinary field theory (much simpler than classical electrodynamics, for example, since there is only *one* field (Ψ), instead of *two* (\mathbf{E} and \mathbf{B}), and it's a *scalar*). It is the bizarre role of the measurement process that gives quantum mechanics its extraordinary richness and subtlety. But what, exactly, *is* a measurement? What makes it so different from other physical processes?¹³ And how can we tell when a measurement has occurred?

Schrödinger posed the essential question most starkly, in his famous **cat paradox**.¹⁴

A cat is placed in a steel chamber, together with the following hellish contraption. . . . In a Geiger counter there is a tiny amount of radioactive substance, so tiny that maybe within an hour one of the atoms decays, but equally probably none of them decays. If one decays then the counter triggers and via a relay activates a little hammer which breaks a container of cyanide. If one has left this entire system for an hour, then one would say the cat is living if no atom has decayed. The first decay would have poisoned it. The wave function of the entire system would express this by containing equal parts of the living and dead cat.

At the end of the hour, the wave function of the cat has the schematic form

$$\psi = \frac{1}{\sqrt{2}}(\psi_{\text{alive}} + \psi_{\text{dead}}). \quad [12.18]$$

¹³There is a school of thought that rejects this distinction, holding that the system and the measuring apparatus should be described by one great big wave function which itself evolves according to the Schrödinger equation. In such theories there is no collapse of the wave function, but one must typically abandon any hope of describing individual events—quantum mechanics (in this view) applies only to *ensembles* of identically prepared systems. See, for example, Philip Pearle *Am. J. Phys.* **35**, 742 (1967), or Leslie E. Ballentine, *Quantum Mechanics: A Modern Development*, 2nd ed., World Scientific, Singapore (1998).

¹⁴E. Schrödinger, *Naturwiss.* **48**, 52 (1935); translation by Josef M. Jauch, *Foundations of Quantum Mechanics*, Addison-Wesley, Reading (1968), p. 185.

The cat is neither alive nor dead, but rather a linear combination of the two, until a measurement occurs—until, say, you peek in the window to check. At that moment your observation forces the cat to “take a stand”: dead or alive. And if you find him to be dead, then it's really *you* who killed him, by looking in the window.

Schrödinger regarded this as patent nonsense, and I think most physicists would agree with him. There is something absurd about the very idea of a *macroscopic* object being in a linear combination of two palpably different states. An electron can be in a linear combination of spin up and spin down, but a cat simply cannot *be* in a linear combination of alive and dead. How are we to reconcile this with the orthodox interpretation of quantum mechanics?

The most widely accepted answer is that the triggering of the Geiger counter constitutes the “measurement,” in the sense of the statistical interpretation, not the intervention of a human observer. It is the essence of a measurement that some *macroscopic* system is affected (the Geiger counter, in this instance). The measurement occurs at the moment when the microscopic system (described by the laws of quantum mechanics) interacts with the macroscopic system (described by the laws of classical mechanics) in such a way as to leave a permanent record. The macroscopic system itself is not permitted to occupy a linear combination of distinct states.¹⁵

I would not pretend that this is an entirely satisfactory resolution, but at least it avoids the stultifying solipsism of Wigner and others, who persuaded themselves that it is the involvement of human consciousness that constitutes a measurement in quantum mechanics. Part of the problem is the word “measurement” itself, which certainly carries a suggestion of human participation. Heisenberg proposed the word “event,” which might be preferable. But I'm afraid “measurement” is so ingrained by now that we're stuck with it. And, in the end, no manipulation of the terminology can completely exorcise this mysterious ghost.

12.5 THE QUANTUM ZENO PARADOX

The collapse of the wave function is undoubtedly the *most* peculiar feature of this whole bizarre story. It was introduced on purely theoretical grounds, to account for the fact that an immediately repeated measurement reproduces the same value. But surely such a radical postulate must carry directly observable consequences. In

¹⁵Of course, in some ultimate sense the macroscopic system is *itself* described by the laws of quantum mechanics. But wave functions, in the first instance, describe individual elementary particles; the wave function of a macroscopic object would be a monstrously complicated composite, built out of all the wave functions of its 10^{23} constituent particles. Presumably somewhere in the statistics of large numbers macroscopic linear combinations become extremely improbable. Indeed, if you *were* able somehow to get a damped pendulum (say) into a linear combination of macroscopically distinct quantum states, it would, in a tiny fraction of the damping time, revert to an ordinary classical state. This phenomenon is called **decoherence**. See, for example, R. Omnes, *The Interpretation of Quantum Mechanics* (Princeton, 1994), Chapter 7.

1977 Misra and Sudarshan¹⁶ proposed what they called the **quantum Zeno effect** as a dramatic experimental demonstration of the collapse of the wave function. Their idea was to take an unstable system (an atom in an excited state, say), and subject it to repeated measurements. Each observation collapses the wave function, resetting the clock, and it is possible by this means to delay indefinitely the expected transition to the lower state.¹⁷

Specifically, suppose a system starts out in the excited state ψ_2 , which has a natural lifetime τ for transition to the ground state ψ_1 . Ordinarily, for times substantially less than τ , the probability of a transition is proportional to t (see Equation 9.42); in fact, since the transition rate is $1/\tau$,

$$P_{2 \rightarrow 1} = \frac{t}{\tau}. \quad [12.19]$$

If we make a measurement after a time t , then, the probability that the system is still in the *upper* state is

$$P_1(t) = 1 - \frac{t}{\tau}. \quad [12.20]$$

Suppose we *do* find it to be in the upper state. In that case the wave function collapses back to ψ_2 , and the process starts all over again. If we make a *second* measurement, at $2t$, the probability that the system is *still* in the upper state is evidently

$$\left(1 - \frac{t}{\tau}\right)^2 \approx 1 - \frac{2t}{\tau}, \quad [12.21]$$

which is the same as it would have been had we never made the first measurement at t . This is what one would naively expect; if it were the whole story there would be nothing gained by repeatedly observing the system, and there would be no quantum Zeno effect.

However, for *extremely* short times, the probability of a transition is *not* proportional to t , but rather to t^2 (see Equation 9.39):¹⁸

$$P_{2 \rightarrow 1} = \alpha t^2. \quad [12.22]$$

In this case the probability that the system is still in the upper state after the two measurements is

$$\left(1 - \alpha t^2\right)^2 \approx 1 - 2\alpha t^2, \quad [12.23]$$

¹⁶B. Misra and E. C. G. Sudarshan, *J. Math. Phys.* **18**, 756 (1977).

¹⁷This effect doesn't have much to do with Zeno, but it *is* reminiscent of the old adage, "a watched pot never boils," so it is sometimes called the **watched pot phenomenon**.

¹⁸In the argument leading to linear time dependence, we assumed that the function $\sin^2(\Omega t/2)/\Omega^2$ in Equation 9.39 was a sharp spike. However, the *width* of the "spike" is of order $\Delta\omega = 4\pi/t$, and for *extremely* short t this approximation fails, and the integral becomes $(t^2/4) \int \rho(\omega) d\omega$.

whereas if we had never made the first measurement it would have been

$$1 - \alpha(2t)^2 \approx 1 - 4\alpha t^2. \quad [12.24]$$

Evidently our observation of the system after time t decreased the net probability of a transition to the lower state!

Indeed, if we examine the system at n regular intervals, from $t = 0$ out to $t = T$ (that is, we make measurements at $T/n, 2T/n, 3T/n, \dots, T$), the probability that the system is still in the upper state at the end is

$$\left(1 - \alpha(T/n)^2\right)^n \approx 1 - \frac{\alpha}{n} T^2, \quad [12.25]$$

which goes to 1 in the limit $n \rightarrow \infty$: A *continuously* observed unstable system never decays at all! Some authors regard this as an absurd conclusion, and a proof that the collapse of the wave function is fallacious. However, their argument hinges on a rather loose interpretation of what constitutes "observation." If the track of a particle in a bubble chamber amounts to "continuous observation," then the case is closed, for such particles certainly do decay (in fact, their lifetime is not measurably extended by the presence of the detector). But such a particle is only intermittently interacting with the atoms in the chamber, and for the quantum Zeno effect to occur the successive measurements must be made *extremely* rapidly, in order to catch the system in the t^2 regime.

As it turns out, the experiment is impractical for spontaneous transitions, but it can be done using *induced* transitions, and the results are in excellent agreement with the theoretical predictions.¹⁹ Unfortunately, this experiment is not as compelling a confirmation of the collapse of the wave function as its designers hoped; the observed effect can be accounted for in other ways.²⁰

In this book I have tried to tell a consistent and coherent story: The wave function (Ψ) represents the state of a particle (or system); particles do not in general possess specific dynamical properties (position, momentum, energy, angular momentum, etc.) until an act of measurement intervenes; the probability of getting a particular value in any given experiment is determined by the statistical interpretation of Ψ ; upon measurement the wave function collapses, so that an immediately repeated measurement is certain to yield the same result. There are other possible interpretations—nonlocal hidden variable theories, the "many worlds" picture, "consistent histories," ensemble models, and others—but I believe this one

¹⁹W. M. Itano, D. J. Heinzen, J. J. Bollinger, and D. J. Wineland, *Phys. Rev. A* **41**, 2295 (1990).

²⁰L. E. Ballentine, *Found. Phys.* **20**, 1329 (1990); T. Petrosky, S. Tasaki, and I. Prigogine, *Phys. Lett. A* **151**, 109 (1990).

is conceptually the *simplest*, and certainly it is the one shared by most physicists today.²¹ It has stood the test of time, and emerged unscathed from every experimental challenge. But I cannot believe this is the end of the story; at the very least, we have much to learn about the nature of measurement and the mechanism of collapse. And it is entirely possible that future generations will look back, from the vantage point of a more sophisticated theory, and wonder how we could have been so gullible.

²¹See Daniel Styer *et al.*, *Am. J. Phys.* **70**, 288 (2002).

APPENDIX

LINEAR ALGEBRA

Linear algebra abstracts and generalizes the arithmetic of ordinary vectors, such as those we encounter in first-year physics. The generalization is in two directions: (1) We allow the scalars to be *complex* numbers, and (2) we do not restrict ourselves to three dimensions.

A.1 VECTORS

A **vector space** consists of a set of **vectors** ($|\alpha\rangle, |\beta\rangle, |\gamma\rangle, \dots$), together with a set of **scalars** (a, b, c, \dots),¹ which is **closed**² under two operations: vector addition and scalar multiplication.

- **Vector Addition**

The “sum” of any two vectors is another vector:

$$|\alpha\rangle + |\beta\rangle = |\gamma\rangle. \quad [\text{A.1}]$$

Vector addition is **commutative**:

$$|\alpha\rangle + |\beta\rangle = |\beta\rangle + |\alpha\rangle, \quad [\text{A.2}]$$

¹For our purposes, the scalars will be ordinary complex numbers. Mathematicians can tell you about vector spaces over more exotic fields, but such objects play no role in quantum mechanics. Note that $\alpha, \beta, \gamma \dots$ are *not* (ordinarily) numbers; they are *names* (labels)—“Charlie,” for instance, or “F43A-9GL,” or whatever you care to use to identify the vector in question.

²That is to say, these operations are always well-defined, and will never carry you outside the vector space.